

**PRACE NAUKOWE**

Uniwersytetu Ekonomicznego we Wrocławiu

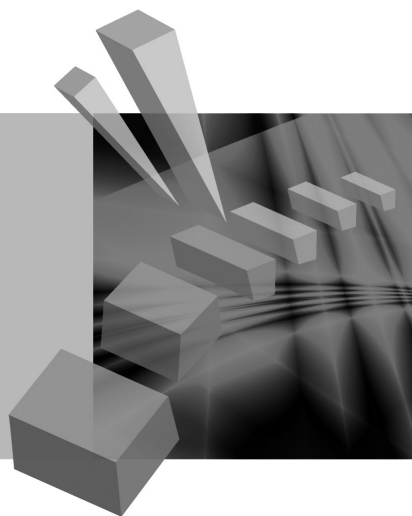
**RESEARCH PAPERS**

of Wrocław University of Economics

**279**

# Taksonomia 21

## Klasyfikacja i analiza danych – teoria i zastosowania



Redaktorzy naukowi

**Krzysztof Jajuga**

**Marek Walesiak**



Wydawnictwo Uniwersytetu Ekonomicznego we Wrocławiu  
Wrocław 2013

Redaktor Wydawnictwa: Aleksandra Śliwka

Redaktor techniczny: Barbara Łopusiewicz

Korektor: Barbara Cibis

Łamanie: Małgorzata Czupryńska

Projekt okładki: Beata Dębska

Publikacja jest dostępna w Internecie na stronach:

[www.ibuk.pl](http://www.ibuk.pl), [www.ebscohost.com](http://www.ebscohost.com),

The Central and Eastern European Online Library [www.ceeol.com](http://www.ceeol.com),

a także w adnotowanej bibliografii zagadnień ekonomicznych BazEkon

[http://kangur.uek.krakow.pl/bazy\\_ae/bazekon/nowy/index.php](http://kangur.uek.krakow.pl/bazy_ae/bazekon/nowy/index.php)

Informacje o naborze artykułów i zasadach recenzowania znajdują się

na stronie internetowej Wydawnictwa

[www.wydawnictwo.ue.wroc.pl](http://www.wydawnictwo.ue.wroc.pl)

Tytuł dofinansowany ze środków Narodowego Banku Polskiego

oraz ze środków Sekcji Klasyfikacji i Analizy danych PTS

Kopiowanie i powielanie w jakiegokolwiek formie

wymaga pisemnej zgody Wydawcy

© Copyright by Uniwersytet Ekonomiczny we Wrocławiu

Wrocław 2013

**ISSN 1899-3192** (Prace Naukowe Uniwersytetu Ekonomicznego we Wrocławiu)

**ISSN 1505-9332** (Taksonomia)

Wersja pierwotna: publikacja drukowana

Druk: Drukarnia TOTEM

## Spis treści

<b>Wstęp</b> .....	9
<b>Sabina Denkowska, Kamil Fijorek, Marcin Salamaga, Andrzej Sokolowski:</b> Sejm VI kadencji – maszynka do głosowania .....	11
<b>Barbara Pawelek, Adam Sagan:</b> Zmienne ukryte w modelach ekonomicznych – respecyfikacja modelu Kleina I .....	19
<b>Jan Paradysz:</b> Nowe możliwości badania koniunktury na rynku pracy .....	29
<b>Krzysztof Najman:</b> Samouczące się sieci GNG w grupowaniu dynamicznym zbiorów o wysokim wymiarze .....	41
<b>Kamila Migdał-Najman:</b> Zastosowanie jednowymiarowej sieci SOM do wyboru cech zmiennych w grupowaniu dynamicznym .....	48
<b>Aleksandra Matuszewska-Janica, Dorota Witkowska:</b> Zróżnicowanie płac ze względu na płeć: zastosowanie drzew klasyfikacyjnych .....	58
<b>Iwona Foryś, Ewa Putek-Szeląg:</b> Przestrzenna klasyfikacja gmin ze względu na sprzedaż użytków gruntowych zbywanych przez ANR w województwie zachodniopomorskim .....	67
<b>Joanna Banaś, Małgorzata Machowska-Szewczyk:</b> Klasyfikacja internetowych rachunków bankowych z uwzględnieniem zmiennych symbolicznych.....	77
<b>Marta Jaročka:</b> Wpływ metody doboru cech diagnostycznych na wynik porządkowania liniowego na przykładzie rankingu polskich uczelni .....	85
<b>Anna Zamojska:</b> Badanie zgodności rankingów wyznaczonych według różnych wskaźników efektywności zarządzania portfelem na przykładzie funduszy inwestycyjnych.....	95
<b>Dorota Rozmus:</b> Porównanie dokładności taksonomicznej metody propagacji podobieństwa oraz zagregowanych algorytmów taksonomicznych opartych na idei metody <i>bagging</i> .....	106
<b>Ewa Wędrowska:</b> Wrażliwość miar dywergencji jako mierników niepodobieństwa struktur.....	115
<b>Katarzyna Wójcik, Janusz Tuchowski:</b> Wpływ automatycznego tłumaczenia na wyniki automatycznej identyfikacji charakteru opinii konsumenckich ...	124
<b>Małgorzata Misztal:</b> Ocena wpływu wybranych metod imputacji na wyniki klasyfikacji obiektów w modelach drzew klasyfikacyjnych.....	135
<b>Anna Czapkiewicz, Beata Basiura:</b> Badanie wpływu wyboru współczynnika zależności na grupowanie szeregów czasowych .....	146
<b>Tomasz Szubert:</b> Czynniki różnicujące poziom zadowolenia z życia oraz wartości życiowe osób sprawnych i niepełnosprawnych w świetle badań „Diagnozy społecznej” .....	154

<b>Marcin Szymkowiak:</b> Konstrukcja estymatorów kalibracyjnych wartości globalnej dla różnych funkcji odległości .....	164
<b>Wojciech Roszka:</b> Szacowanie łącznych charakterystyk cech nieobserwowanych łącznie .....	174
<b>Justyna Brzezińska:</b> Metody wizualizacji danych jakościowych w programie <b>R</b> .....	182
<b>Agata Sielska:</b> Regionalne zróżnicowanie potencjału konkurencyjnego polskich gospodarstw rolnych w województwach po akcesji do Unii Europejskiej .....	191
<b>Mariusz Kubus:</b> Liniowy model prawdopodobieństwa z regularyzacją jako metoda doboru zmiennych .....	201
<b>Beata Basiura:</b> Metoda Warda w zastosowaniu klasyfikacji województw Polski z różnymi miarami odległości .....	209
<b>Katarzyna Wardzińska:</b> Wykorzystanie metody obwiedni danych w procesie klasyfikacji przedsiębiorstw .....	217
<b>Katarzyna Dębowska:</b> Modelowanie upadłości przedsiębiorstw oparte na próbach niezbilansowanych .....	226
<b>Danuta Tarka:</b> Wpływ metody doboru cech diagnostycznych na wyniki klasyfikacji obiektów na przykładzie danych dotyczących ochrony środowiska ..	235
<b>Artur Czech:</b> Zastosowanie wybranych metod doboru zmiennych diagnostycznych w badaniach konsumpcji w ujęciu pośrednim .....	246
<b>Beata Bal-Domańska:</b> Ocena relacji zachodzących między inteligentnym rozwojem a spójnością ekonomiczną w wymiarze regionalnym z wykorzystaniem modeli panelowych .....	255
<b>Mariola Chrzanowska:</b> <i>Ordinary kriging</i> i <i>inverse distance weighting</i> jako metody szacowania cen nieruchomości na przykładzie warszawskiego rynku .....	264
<b>Adam Depta:</b> Zastosowanie analizy wariancji w badaniu jakości życia na podstawie kwestionariusza SF-36v2 .....	272
<b>Maciej Beręsewicz, Tomasz Klimanek:</b> Wykorzystanie estymacji pośredniej uwzględniającej korelację przestrzenną w badaniach cen mieszkań .....	281
<b>Karolina Paradysz:</b> Benchmarkowa analiza estymacji dla małych obszarów na lokalnych rynkach pracy .....	291
<b>Anna Gryko-Nikitin:</b> Dobór parametrów w równoległych algorytmach genetycznych dla problemu plecakowego .....	301
<b>Tomasz Ząbkowski, Piotr Jałowiecki:</b> Zastosowanie reguł asocjacyjnych do analizy danych ankietowych w wybranych obszarach logistyki przedsiębiorstw przetwórstwa rolno-spożywczego .....	311
<b>Agnieszka Przedborska, Małgorzata Misztal:</b> Zastosowanie metod statystyki wielowymiarowej do oceny wydolności stawów kolanowych u pacjentów z chorobą zwyrodnieniową leczonych operacyjnie .....	321
<b>Dorota Perło:</b> Rozwój zrównoważony w wymiarze gospodarczym, społecznym i środowiskowym – analiza przestrzenna .....	331

<b>Ewa Putek-Szeląg, Urszula Gieraltowska, Analiza i diagnoza wielkości produkcji energii odnawialnej w Polsce na tle krajów Unii Europejskiej..</b>	342
--	-----

## Summaries

<b>Sabina Denkowska, Kamil Fijorek, Marcin Salamaga, Andrzej Sokolowski: VIth-term Sejm – a voting machine .....</b>	18
<b>Barbara Pawelek, Adam Sagan: Latent variables in econometric models – respecification of Klein I model .....</b>	28
<b>Jan Paradysz: New possibilities for studying the situation on the labour market .....</b>	40
<b>Krzysztof Najman: Self-learning neural network of GNG type in the dynamic clustering of high-dimensional data.....</b>	47
<b>Kamila Migdał-Najman: Applying the one-dimensional SOM network to select variables in dynamic clustering .....</b>	57
<b>Aleksandra Matuszewska-Janica, Dorota Witkowska: Gender wage gap: application of classification trees.....</b>	66
<b>Iwona Foryś, Ewa Putek-Szeląg: Spatial classification of communes by usable land traded by the APA in the Zachodniopomorskie voivodeship...</b>	76
<b>Joanna Banaś, Małgorzata Machowska-Szewczyk: Classification of Internet banking accounts including symbolic variables .....</b>	84
<b>Marta Jarocka: The impact of the method of the selection of diagnostic variables on the result of linear ordering on the example of ranking of universities in Poland.....</b>	94
<b>Anna Zamojska: Empirical analysis of the consistency of mutual fund ranking for different portfolio performance measures.....</b>	105
<b>Dorota Rozmus: Comparison of accuracy of affinity propagation clustering and cluster ensembles based on bagging idea.....</b>	114
<b>Ewa Wędrowska: Sensitivity of divergence measures as structure dissimilarity measurements .....</b>	123
<b>Katarzyna Wójcik, Janusz Tuchowski: Machine translation impact on the results of the sentiment analysis .....</b>	134
<b>Małgorzata Misztal: Assessment of the influence of selected imputation methods on the results of object classification using classification trees ...</b>	145
<b>Anna Czapkiewicz, Beata Basiura: Simulation study of the selection of coefficient depending on the clustering time series.....</b>	153
<b>Tomasz Szubert: Factors differentiating the level of satisfaction with life and the life's values of people with and without disabilities in the light of the "Social Diagnosis" survey .....</b>	162
<b>Marcin Szymkowiak: Construction of calibration estimators of totals for different distance measures .....</b>	173

<b>Wojciech Roszka:</b> Joint characteristics' estimation of variables not jointly observed.....	181
<b>Justyna Brzezińska:</b> Visualizing categorical data in $\mathbf{R}$ .....	190
<b>Agata Sielska:</b> Regional diversity of competitiveness potential of Polish farms after the accession to the European Union .....	200
<b>Mariusz Kubus:</b> Regularized linear probability model as a filter .....	208
<b>Beata Basiura:</b> The Ward method in the application for classification of Polish voivodeships with different distances.....	216
<b>Katarzyna Wardzińska:</b> Application of Data Envelopment Analysis in company classification process.....	225
<b>Katarzyna Dębowska:</b> Modeling corporate bankruptcy based on unbalanced samples .....	234
<b>Danuta Tarka:</b> Influence of the features selection method on the results of objects classification using environmental data.....	245
<b>Artur Czech:</b> Application of chosen methods for the selection of diagnostic variables in indirect consumption research.....	254
<b>Beata Bal-Domańska:</b> Assessment of relations occurring between smart growth and economic cohesion in regional dimension using panel models .....	263
<b>Mariola Chrzanowska:</b> Ordinary kriging and inverse distance weighting as methods of estimating prices based on Warsaw real estate market .....	271
<b>Adam Depta:</b> Application of analysis of variance in the study of the quality of life based on questionnaire SF-36v2 .....	280
<b>Maciej Beręsewicz, Tomasz Klimanek:</b> Using indirect estimation with spatial autocorrelation in dwelling price surveys.....	290
<b>Karolina Paradysz:</b> Benchmark analysis of small area estimation on local labor markets .....	300
<b>Anna Gryko-Nikitin:</b> Selection of various parameters of parallel evolutionary algorithm for knapsack problems .....	310
<b>Tomasz Ząbkowski, Piotr Jałowiecki:</b> Application of association rules for the survey of data analysis in the selected areas of logistics in food processing companies .....	320
<b>Agnieszka Przedborska, Małgorzata Misztal:</b> Using multivariate statistical methods to assess the capacity of the knee joint among the patients treated surgically for osteoarthritis .....	330
<b>Dorota Perło:</b> Sustainable development in the economic, social and environmental dimensions – spatial analysis.....	341
<b>Ewa Putek-Szeląg, Urszula Gieraltowska:</b> Analysis and diagnosis of the volume of renewable energy production in Poland compared to EU countries .....	352

**Maciej Beręsewicz, Tomasz Klimanek**

Uniwersytet Ekonomiczny w Poznaniu

---

## **WYKORZYSTANIE ESTYMACJI POŚREDNIEJ UWZGLĘDNIAJĄCEJ KORELACJĘ PRZESTRZENNĄ W BADANIACH CEN MIESZKAŃ**

---

**Streszczenie:** Artykuł przedstawia propozycję wykorzystania metod estymacji pośredniej (w tym także metody, która uwzględnia korelację przestrzenną) do oszacowania pewnych charakterystyk rynku nieruchomości w województwie wielkopolskim. W konstrukcji odpowiednich estymatorów statystyki małych obszarów autorzy postanowili wykorzystać, oprócz modeli przekrojowych, także najnowsze podejścia w estymacji pośredniej wykorzystujące zależności przestrzenne. Ze względu na utrudniony dostęp do danych transakcyjnych badania oparte zostały na danych ofertowych zawierających informację o lokalizacji nieruchomości w przestrzeni geograficznej (dane zorientowane przestrzennie).

**Słowa kluczowe:** statystyka małych obszarów, autokorelacja przestrzenna, analiza rynku nieruchomości.

### **1. Wstęp**

W ostatnich latach obserwuje się wzrost zainteresowania metodami estymacji pośredniej w Polsce. Wiele z dotychczasowych zastosowań dotyczyło problemów zwiększenia pokrycia informacyjnego dla potrzeb statystyki publicznej [Dehnel, Gołata 2006; Gołata 2004; Klimanek, Szymkowiak 2012; Kubacki 2008; Paradysz 2003]. Jednakże od samego początku stosowania metod statystyki małych obszarów podkreślano, że naturalnym odbiorcą wyników szacunków wydaje się obok statystyki publicznej i instytucji samorządowych także szeroko rozumiany biznes.

Należy także zwrócić uwagę na fakt, że większość zastosowań statystyki małych obszarów dotyczy estymacji charakterystyk rynku pracy, demografii i statystyki społecznej; stosunkowo niewielka jest liczba artykułów poświęconych zastosowaniom estymacji pośredniej w badaniach biznesowych.

Celem artykułu jest zastosowanie estymacji pośredniej do oszacowania przeciętnej ceny metra kwadratowego mieszkań na rynku nieruchomości mieszkaniowych w Poznaniu w sierpniu 2012 r.. Zastosowane podejście może zdaniem autorów, stanowić alternatywę w przypadku braku dostępu do danych transakcyjnych. Ponadto, przyjmując założenie o występowaniu wpływu lokalizacji w mieście na średnią cenę

mieszkania, autorzy podjęli próbę zastosowania w konstrukcji estymatora modelu wykorzystującego autokorelację przestrzenną. Tego rodzaju podejście zostało opisane w projekcie EURAREA, ale liczba zastosowań modeli wciąż jest niewystarczająca i wynika głównie z niedostatecznego wykorzystania istniejących informacji geoprzestrzennych – informacji opartych na współrzędnych geograficznych obiektów lub pewnych szczególnych charakterystykach związanych z tymi punktami.

## 2. Opis procedury badawczej

Rynek nieruchomości charakteryzuje się utrudnionym dostępem do danych charakteryzujących zawierane transakcje kupna-sprzedaży. W Internecie natomiast znajduje się spora liczba portali zajmujących się przedstawianiem ofert dotyczących m.in. sprzedaży mieszkań na rynku zarówno pierwotnym, jak i wtórnym. Stanowią one cenne źródło informacji na temat mieszkań oferowanych przez osoby prywatne oraz pośredników. Portale zawierają ceny ofertowe, co oznacza, że mogą się różnić od cen transakcyjnych, jednak m.in. Narodowy Bank Polski publikuje ceny ofertowe jako element cyklicznych raportów dotyczących rynku nieruchomości<sup>1</sup>, a także tworzy bazę rynku nieruchomości (BaRn), w której uwzględnia zmiany cen zarówno transakcyjnych, jak i ofertowych. Portale, takie jak Gratka, Domy.pl (przy współpracy Open Finance) czy OtoDom, na podstawie cen ofertowych tworzą indeksy cen mieszkań oraz analizują sytuację na rynku nieruchomości<sup>2</sup>.

W związku z ograniczonymi możliwościami uzyskania informacji na temat cen transakcyjnych oraz wykorzystywaniem cen ofertowych w analizach GUS i NBP postanowiono przeprowadzić analizę z wykorzystaniem cen oferowanych na portalach internetowych. W tym celu w programie R z wykorzystaniem pakietów XML oraz RCurl został napisany program (potocznie nazywany „pajakiem internetowym”) umożliwiający automatyczne pobieranie informacji o ofertach mieszkań z rynku pierwotnego i wtórnego. Działanie programu opiera się na następujących krokach:

1. Wejść na stronę wyników mieszkań dla Poznania. Ustal  $i = 1$  oraz  $n$ , które oznacza ostatnią stronę.

1.1 Pobierz ze strony wszystkie linki, które dotyczą ofert mieszkań ( $j = 1 \dots m$ ).

1.2 Dla każdego linku ( $j = 1 \dots m$ ) z punktu 1.1 wejdź na stronę i pobierz informacje o mieszkaniu.

1.3 Jeżeli pobrano informacje o wszystkich mieszkaniach z 1.1, wróć do punktu 1.

---

<sup>1</sup> Raporty można znaleźć na stronie: [http://www.nbp.pl/home.aspx?f=/publikacje/rynek\\_nieruchomosci/index1.html](http://www.nbp.pl/home.aspx?f=/publikacje/rynek_nieruchomosci/index1.html).

<sup>2</sup> Opracowania i raporty można znaleźć na stronach internetowych [www.gratka.pl](http://www.gratka.pl), [www.domy.pl](http://www.domy.pl) czy [www.otodom.pl](http://www.otodom.pl).



2. Przejdź na kolejną stronę wyszukiwań ( $\mathbf{i} = \mathbf{i} + \mathbf{1}$ ). Jeżeli nie jest to ostatnia strona ( $\mathbf{i} \neq \mathbf{n}$ ), wróć do punktu 1, w przeciwnym wypadku przejdź do punktu 3.

3. Zakończ działanie pętli.

Program działa do momentu, aż odwiedzi wszystkie strony zawierające wyniki wyszukiwań zadanych na początku (Poznań, rynek wtórny) oraz podstrony zawierające informacje o mieszkaniach. Oferty mieszkań były opisane szeregiem zmiennych, m.in. takich jak ceny mieszkania, ceny metra kwadratowego, liczba pokoi, powierzchnia, jak również informacja na temat położenia (w postaci współrzędnych geograficznych). Osoby zajmujące się publikowaniem informacji o danej nieruchomości umieszczały dodatkowe informacje, które nie występowały dla wszystkich ofert (np. typ budynku, rok budowy, stan).

Na potrzeby artykułu z portalu Domy.pl pobrano 14 229 ofert dotyczących mieszkań z wtórnego rynku nieruchomości mieszkaniowych w Poznaniu w sierpniu 2012 r., które po zastosowaniu omówionych w kolejnym rozdziale technik czyszczenia danych ograniczono do 9952 ofert.

Pierwszym etapem czyszczenia uzyskanego zbioru danych była eliminacja obserwacji powtarzających się. Deduplikacji dokonano, porównując ID oraz linki do ofert mieszkań. Następnie usunięto z analizy oferty mieszkań, które nie zawierały informacji o powierzchni, liczbie pokoi oraz położeniu.

Kolejnym etapem było wykorzystanie metody Least Trimmed Squares znajdującej się w poleceniu *PROC ROBUSTREG* pakietu SAS w celu detekcji wartości odstających. Metoda ta została zaproponowana przez Rousseeuwa [1984] i ma następującą postać:

$$\hat{\theta}_{LTS} = \arg \min_{\theta} Q_{LTS}(\theta), \quad (1)$$

gdzie:

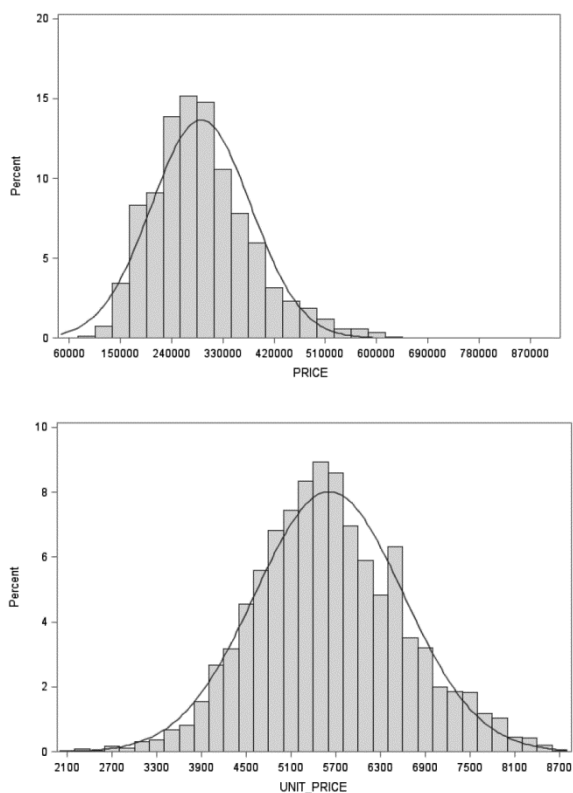
$$Q_{LTS}(\theta) = \sum_{i=1}^h r_i^2, \quad (2)$$

$r_{(1)}^2 \leq r_{(2)}^2 \leq \dots \leq r_{(n)}^2$  są uszeregowanymi kwadratami reszt  $r_i^2 = (y_i - x_i^T \theta)^2$ ,  $i = 1, \dots, n$ , a  $h$  określone jest w przedziale  $\frac{n}{2} + 1 \leq h \leq \frac{3n+p+1}{4}$ .

W pakiecie SAS wartość progowa jest domyślnie ustawiona na  $h = \frac{3n+p+1}{4}$ .

Efektom etapu czyszczenia danych była baza, którą na potrzeby badania nazwano pseudopopulacją mieszkań w Poznaniu w sierpniu 2012 r. Na rysunku 1 zaprezentowano rozkład cen mieszkań oraz cen metra kwadratowego mieszkań w pseudopopulacji, tzn. po zastosowaniu procedur czyszczenia danych. Rozkład ceny metra kwadratowego jest bardziej zbliżony do rozkładu normalnego niż rozkład cen mieszkań.

W związku z faktem, że celem badania symulacyjnego miało być oszacowanie średniej ceny metra kwadratowego według jednostek ewidencyjnych miasta Pozna-

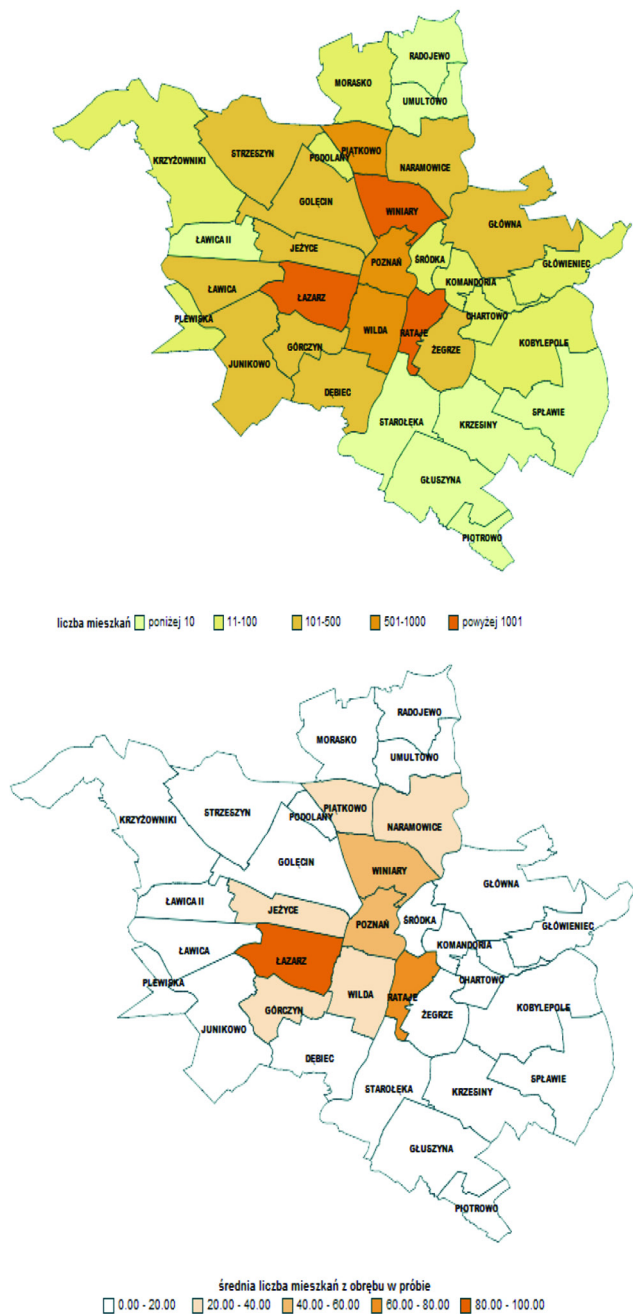


**Rys. 1.** Rozkład ceny oraz (PRICE) ceny m<sup>2</sup> (UNIT\_PRICE) mieszkań w pseudopopulacji

Źródło: opracowanie własne w pakiecie SAS.

nia (obrębów<sup>3</sup>), a w populacji zawarta była jedynie informacja o współrzędnych geograficznych mieszkań znajdujących się na rynku ofert, należało przyporządkować mieszkańom identyfikator obrębu. W tym celu posłużono się procedurą *PROC GINSIDE*, która sprawdza, czy podane współrzędne geograficzne mieszkania należą do określonego poligonu (obrębu miasta Poznania), a jeśli tak, to rekord w bazie zostaje uzupełniony o identyfikator tego poligonu. Następnie z populacji tej wylosowano 1000 prób o wielkości 5%, stosując schemat losowania prostego. Zastosowany schemat losowania spowodował, że w przypadku obrębów o niewielkiej liczbie mieszkań w populacji, w wylosowanych próbach reprezentacja mieszkań była niewielka lub nawet równa 0. Rozkład przestrzenny liczebności mieszkań w pseudopopulacji i badaniu symulacyjnym przedstawia rys. 2.

<sup>3</sup> Autorzy zdają sobie sprawę z faktu, że obręb geodezyjny nie jest najlepszym kryterium delimitacji przestrzeni dla wyznaczenia obszarów cenności, gdyż te przekraczają granice administracyjne (geodezyjne). Przyjęcie takiej jednostki przestrzennej wynikało z określonego zdefiniowania małego obszaru (domeny), dla której można było uzyskać podkład dla mapy numerycznej.



**Rys. 2.** Przestrzenny rozkład liczby mieszkań w pseudopopulacji i badaniu symulacyjnym  
 Źródło: opracowanie własne w pakiecie SAS.

Dla każdej z prób do oszacowania średniej ceny metra kwadratowego w obrębach miasta Poznania zastosowano 4 estymatory, przy czym w przypadku estymatorów typu GREG, EBLUP\_B i SEBLUP, które wykorzystują zmienne pomocnicze z próby i spoza próby, do modelowania ceny metra kwadratowego wykorzystano powierzchnię mieszkania i liczbę pokoi. Należy zwrócić uwagę, że jakość modelu jest słaba, gdyż zmienność ceny metra kwadratowego została wyjaśniona przez wybrane zmienne objaśniające jedynie w około 12%.

Charakterystykę modelu w pseudopopulacji przedstawia tab. 1.

**Tabela 1.** Charakterystyka modelu w pseudopopulacji

Zmienna	Liczba stopni swobody	Oszacowanie parametru	Błąd standardowy	Statystyka t-Studenta	Prawdopodobieństwo testowe
Wyraz wolny	1	6745,65	32,58	207,05	<,0001
AREA	1	-16,00	0,98	-16,36	<,0001
ROOMS	1	-135,55	18,27	-7,42	<,0001

$$R^2 = 0,123, \text{ współczynnik zmienności losowej} = 16,2\%$$

Źródło: obliczenia własne w pakiecie SAS.

Zastosowano estymatory<sup>4</sup>:

- estymator bezpośredni (Horvitz-Thompsona)

$$\hat{Y}_d^{DIRECT} = \frac{1}{\hat{N}_d} \sum_{i \in u_d} w_{id} y_{id}, \quad (3)$$

gdzie  $\hat{N}_d = \sum_{i \in u_d} w_{id}$  oraz  $w_{id} = \frac{1}{\pi_{id}}$

przy założeniu, że  $\pi_{id,jd} = 0$  dla wszystkich  $d \neq d'$  lub  $i \neq j$ ;

- estymator GREG

$$y_{id} = \mathbf{x}_{id}^T \boldsymbol{\beta} + \varepsilon_{id}, \quad (4)$$

gdzie  $E(\varepsilon_{id}) = 0$ ,  $Var(\varepsilon_{id}) = \sigma_\varepsilon^2$ ,

$$\hat{Y}_d^{GREG} = \frac{1}{\hat{N}_d} \sum_{i \in s_d} \frac{y_i}{\pi_i} + \left( \bar{\mathbf{x}}_d^T - \frac{1}{\hat{N}_d} \sum_{i \in s_d} \frac{\mathbf{x}_i}{\pi_i} \right)^T \hat{\boldsymbol{\beta}}, \quad (5)$$

gdzie  $\hat{N}_d = \sum_{i \in s_d} \frac{1}{\pi_i}$  i  $\hat{\boldsymbol{\beta}}$  są oszacowane z wykorzystaniem ważonej metody najmniejszych kwadratów poprzez użycie wag wynikających ze schematu losowania:

<sup>4</sup> Wzory na oszacowania błędów średniokwadratowych zostały pominięte ze względu na ograniczenia objętości tekstu niniejszej publikacji. Są one umieszczone w dokumentacji projektu EURAREA na stronie Urzędu Statystycznego Wielkiej Brytanii – <http://www.statistics.gov.uk/eurarea>.

$$\hat{\beta} = \left( \sum_{i \in u_d} w_{id} x_{id} x_{id}^T \right)^{-1} \sum_{i \in u_d} w_{id} x_{id} y_{id} ; \quad (6)$$

- estymator EBLUP\_B będący kombinacją liniową estymatora bezpośredniego i syntetycznego (*EURAREA\_Project\_Reference\_Volume* 2004),

$$\hat{Y}_d^{EBLUP\_B} = \gamma_d \hat{Y}_d^{DIRECT} + (1 - \gamma_d) \bar{X}_d^T \hat{\beta} \quad (7)$$

$$\gamma_d = \frac{\hat{\sigma}_u^2}{\hat{\sigma}_u^2 + \hat{\sigma}_e^2}, \text{ przy czym } u_d \sim iid N(0, \sigma_u^2), \quad e_{id} \sim iid N(0, \sigma_e^2)$$

$$\hat{\beta} = (x^T D^{-1} x)^{-1} x^T D^{-1} y,$$

gdzie:  $y$  – wektor obserwacji na zmiennej objaśnianej,

$x$  – macierz o wierszach składających się z  $\bar{x}_d^T$ ,

$D$  – macierz o iteracyjnie aktualizowanych elementach  $(\hat{\sigma}_u^2 + \hat{\sigma}_e^2)$  na diagonalu;

- estymator SEBLUP<sup>5</sup> uwzględniający autokorelację efektów losowych związanych z lokalizacją domen w przestrzeni [Saei, Chambers 2004; D'Alò, Falorsi, Solari 2004].

W zapisie macierzowym model można zapisać następująco:

$$y = X\beta + Zu + e, \quad (8)$$

gdzie:  $y$  jest wektorem zmiennej objaśnianej,  $X$  i  $Z$  są znanymi macierzami rzędu odpowiednio:  $N \times P$  (liczba obserwacji razy liczba zmiennych pomocniczych) i  $N \times D$  (liczba obserwacji razy liczba małych obszarów). Macierz  $Z$  jest macierzą incydencji zdefiniowaną następująco:

$$Z = \begin{bmatrix} 1_{N_a} & \cdot & \cdot & 0 \\ 0 & \cdot & \cdot & 0 \\ 0 & \cdot & \cdot & 0 \\ 0 & \cdot & \cdot & 1_{N_b} \end{bmatrix}, \text{ gdzie } 1_{N_a} \text{ jest wektorem o wymiarach } N_a, \text{ którego}$$

wszystkie elementy są równe 1,  $u$  oraz  $e$  są wektorami zmiennych losowych o wartościach oczekiwanych równych 0 oraz macierzy wariancji – kowariancji odpowiednio:

$$N \sim [0, \sigma_u^2 A] \text{ oraz } N \sim [0, \sigma^2 I_N],$$

elementy  $a_{(dd')}$  macierzy  $A$  są dane wzorem:

<sup>5</sup> SEBLUP – Spatial EBLUP.

$$\mathbf{a}_{(dd')} = \left[ 1 + \delta_{(dd')} \exp\left(\frac{\text{dist}(dd')}{\alpha}\right) \right]^{-1}, \quad (9)$$

gdzie:  $\text{dist}(dd')$  oznacza odległość między małymi obszarami  $d$  i  $d'$ .

$$\delta_{(dd')} = \begin{cases} 0 & \text{for } d = d' \\ 1 & \text{for } d \neq d' \end{cases}, \quad (10)$$

a  $\alpha$  jest parametrem skali.

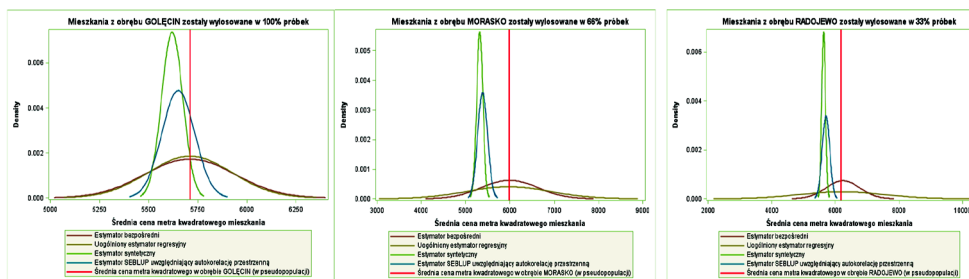
### 3. Uzyskane wyniki

Ze względu na ograniczania objętości tekstu niniejszej publikacji wyniki zostaną przedstawione w sposób bardzo syntetyczny<sup>6</sup>. Jednostki ewidencyjne (obrębny) zostały podzielone na trzy kategorie:

Kategoria A – liczebność próby we wszystkich symulacjach przekraczała 3 mieszkania.

Kategoria B – liczebność próby w ponad 50% symulacji przekraczała 3. Kategoria C – liczebność próby w więcej niż 50% symulacji była mniejsza bądź równa 3.

Dla wymienionych powyżej kategorii wybrano po jednym przykładzie obrębny i dokonano oceny obciążenia 4 zastosowanych estymatorów i oceny wzrokowej wariancji tych estymatorów. Podkreślić jednak należy, że charakterystyki rozkładów są podobne dla pozostałych obrębny w ramach danej kategorii, stąd w ocenie autorów zaprezentowane przypadki są dobrymi reprezentantami obrębny w poszczególnych kategoriach.



Rys. 3. Rozkład estymatorów w obrębie Golecína (kat. A), Moraska (kat. B) oraz Radojewa (kat. C)

Źródło: opracowanie własne w pakiecie SAS.

<sup>6</sup> Autorzy dysponują rozkładem estymatorów dla wszystkich 33 analizowanych obrębny, ale ograniczenia publikacji spowodowały to, że możliwa jest prezentacja jedynie wybranych jednostek ewidencyjnych.

## 4. Wnioski

Przeprowadzone badanie symulacyjne dostarczyło pewnych wniosków, które można sformułować następująco:

- estymator bezpośredni, chociaż nieobciążony, charakteryzuje się dwiema podstawowymi wadami w przypadku estymacji dla małych domen:
  - 1) ma nieakceptowalną wariancję, a w związku z tym także błąd szacunku,
  - 2) w przypadku zerowej próby w domenie nie można wyznaczyć oceny estymatora;
- uogólniony estymator regresyjny, chociaż umożliwia uzyskanie oceny estymatora w przypadku zerowych prób, to jednak charakteryzuje się równie dużą wariancją co estymator bezpośredni;
- estymatory syntetyczne i uwzględniające autokorelację przestrzenną charakteryzują się niewielką wariancją. W porównaniu do estymatorów bezpośrednich są one jednak obciążone;
- analiza przestrzennego rozkładu estymatora uwzględniającego autokorelację przestrzenną sugeruje, że może być on dobrym narzędziem do szacowania cen nieruchomości w domenach, którymi są na przykład części miast o niewielkiej liczbie mieszkań na wtórnym rynku nieruchomości mieszkaniowych.

## Literatura

- aiM Property Małeccy Adamiczka spółka jawna, *Analiza cen transakcyjnych lokali mieszkalnych*, Poznań 2009.
- D'Alò M., Falorsi S., Solari F., *EURAREA Documentation on SAS/IML program on Linear Mixed Model with Spatial Correlated Area Effects in Small Area Estimation*, EURAREA Deliverable, <http://www.ons.gov.uk/ons/guide-method/method-quality/general-methodology/spatial-analysis-and-modelling/eurarea/downloads/index.html>, 2004
- Dehnel G., Gołata E., *Attempts to estimate basic information for small business in Poland*, "Statistics in Transition", Główny Urząd Statystyczny, Warszawa 2006, vol. 6, Number 5, s. 755-776
- EURAREA Project Reference Volume*, <http://www.statistics.gov.uk/eurarea>, 2004.
- Gołata E., *Problems of estimate unemployment for small domains in Poland*, "Statistics in Transition", Główny Urząd Statystyczny, Warszawa 2004, vol. 6, Number 5, s. 755-776.
- Institut Ekonomiczny Narodowego Banku Polskiego, *Raport o sytuacji na rynku nieruchomości mieszkaniowych i komercyjnych w Polsce w 2011 r.*, NBP, Warszawa 2012
- Klimanek T., Szymkowiak M., *Zastosowanie estymacji pośredniej uwzględniającej korelację przestrzenną w opisie niektórych charakterystyk rynku pracy*, [w:] *Taksonomia 19, Klasyfikacja i analiza danych – teoria i zastosowania*, Wydawnictwo Uniwersytetu Ekonomicznego we Wrocławiu, 2012, s. 601-609.
- Kubacki, J., *Application of Bayesian estimation methods for small domains in the Polish Labor Force Survey*, *Acta Universitatis Lodzianis, Folia Oeconomica* 216, Łódź 2008, s. 389-396
- Paradysz J., *Zasilanie publicznej statystyki regionalnej za pomocą estymacji dla małych obszarów w perspektywie wykorzystania rejestrów administracyjnych*, „Wiadomości Statystyczne”, Główny Urząd Statystyczny, Warszawa 2003, nr 4, s. 1-9.
- Rao J.N.K., *Small Area Estimation*, John Wiley & Sons, Inc, 2004.

Rousseeuw P.J., *Least median of squares regression*, "Journal of the American Statistical Association" 1984, 79, 871-880.

Saei A., Chambers R., *Small Area Estimation Under Linear and Generalized Linear Mixed Models With Time and Area Effects*, University of Southampton, 2004.

## **USING INDIRECT ESTIMATION WITH SPATIAL AUTOCORRELATION IN DWELLING PRICE SURVEYS**

**Summary:** The article presents the application of indirect estimation methods (including the method accounting for spatial correlation) to estimate some characteristics of real estate market in Wielkopolska Voivodeship. To build the small area estimators the authors decided to apply not only cross-sectional models but also the most up to date approach using spatial correlation. Because the access to the transactional data was not possible the research was based on the Internet data (offers) which included information about localization (spatially oriented data).

**Keywords:** small area statistics, spatial autocorrelation, real estate market analysis.